

# One-Class Slab Support Vector Machine

Victor Fragoso<sup>†</sup>

Walter Scheirer<sup>‡</sup>

Joao Hespanha<sup>‡</sup>

Matthew Turk<sup>‡</sup>

<sup>†</sup>West Virginia University

<sup>‡</sup>University of Notre Dame

<sup>‡</sup>University of California, Santa Barbara

victor.fragoso@mail.wvu.edu

wscheire@nd.edu

{hespanha@ece, mturk@cs}.ucsb.edu

**Abstract**—This work introduces the one-class slab SVM (OCSSVM), a one-class classifier that aims at improving the performance of the one-class SVM. The proposed strategy reduces the false positive rate and increases the accuracy of detecting instances from novel classes. To this end, it uses two parallel hyperplanes to learn the normal region of the decision scores of the target class. OCSSVM extends one-class SVM since it can scale and learn non-linear decision functions via kernel methods. The experiments on two publicly available datasets show that OCSSVM can consistently outperform the one-class SVM and perform comparable to or better than other state-of-the-art one-class classifiers.

## I. INTRODUCTION

Current recognition systems perform well when their training phase uses a vast amount of samples from all classes encountered at test time. However, these systems significantly decrease in performance when they face the open-set recognition problem [20]: recognition in the presence of samples from unknown or novel classes. This occurs even for already solved datasets (*e.g.*, the Letter dataset [10]) that are recontextualized as open-set recognition problems. The top of the Figure 1 illustrates the general open-set recognition problem.

Recent work has aimed at increasing the robustness of classifiers in this context [1], [19], [20]. However, these approaches assume knowledge of at least a few classes during the training phase. Unfortunately, many recognition systems only have a few samples from just the target class. For example, collecting images from the normal state of a retina is easier than collecting those from abnormal retinas [25].

One-class classifiers are useful in applications where collecting samples from negative classes is challenging, but gathering instances from a target class is easy. An ensemble of one-class classifiers can solve the open-set recognition problem. This is because each one-class classifier can recognize samples of the class it was trained for and detect novel samples; see Figure 1 for an illustration of the ensemble of one-class classifiers. Unlike other solutions to the open-set recognition problem (*e.g.*, the 1-vs-Set SVM [20]), the ensemble offers parallelizable training and easy integration of new categories. These computational advantages follow from the independence of each classifier and allow the ensemble to scale well with the number of target classes.

However, the one-class classification problem is a challenging binary categorization task. This is because the classifier is trained with only positive examples from the target class, yet, it must be able to detect novel samples (negative class data). For instance, a one-class classifier trained to detect normal retinas must learn properties from them to recognize

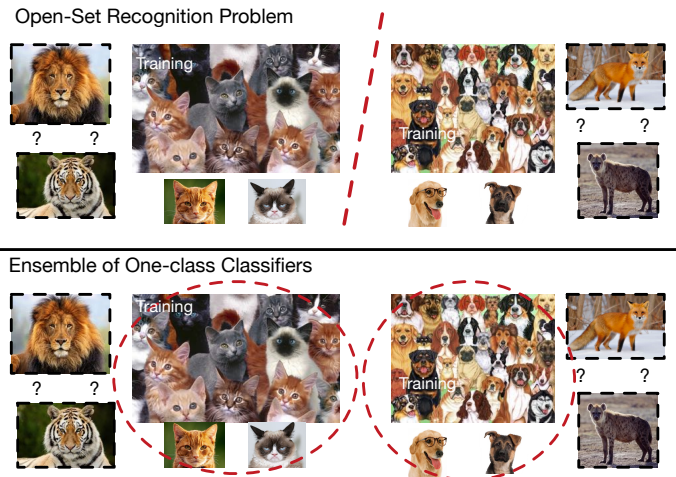


Fig. 1. The open-set recognition problem (top) challenges existing recognition systems. This is because classifiers can face instances from novel or unknown classes (images with dashed-frames). These novel classes cause failures during prediction time. Collecting instances from all the possible classes is a challenging task in many applications. For instance, collecting and labeling instances of all existing animals to avoid this problem is impractical. An ideal solution to this open-set recognition problem is an ensemble of one-class classifiers (bottom). A single one-class classifier only requires instances of a target positive class to train (illustrated as circles). Such classifiers detect samples from the target classes and identify unknown instances. However, their performance needs improvement in order to solve the open-set recognition problem. The proposed approach improves the performance of the one-class SVM. It is a step towards the solution of the open-set recognition problem with an ensemble of one-class classifiers.

other images of normal and abnormal retinas. A vast amount of research has focused on tackling the challenges faced in the one-class classification problem. These strategies include statistical methods [6], [18], neural networks [2], [15], and kernel methods [13], [22], [23].

Despite the advancements, the performance of one-class classifiers falls short for open-set recognition problems. To improve the performance of one-class classifiers, we propose a new algorithm called the one-class slab SVM (OCSSVM), which reduces the rate of classifying instances from a novel class as positive (false positive rate) and increases the rate of detecting instances from a novel class (true negative rate). This work focuses on the one-class SVM classifier as a basis because it can scale well and can learn non-linear decision functions via kernel methods.

The one-class SVM (OCSVM) learns a hyperplane that keeps most of the instances of the target class on its positive side. However, instances from the negative class can also be on the positive side of this hyperplane. The OCSVM does not account for this case, which makes it prone to a high false

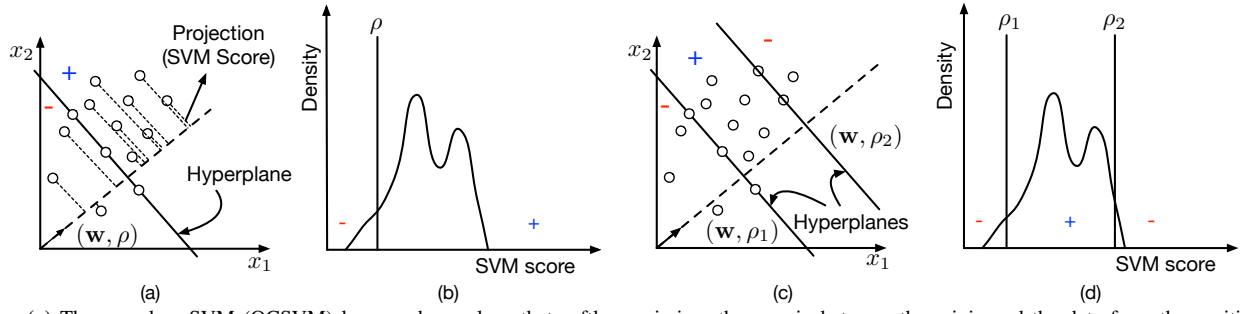


Fig. 2. (a) The one-class SVM (OCSVM) learns a hyperplane that softly maximizes the margin between the origin and the data from the positive class. Its decision function projects the data onto the normal vector  $\mathbf{w}$  to produce the SVM scores. (b) Subsequently, the decision function labels the samples as negative when the SVM scores fall below a threshold  $\rho$ , or labels them as positive otherwise. However, the one-class SVM does not account for outliers that can occur on the right tail of the SVM score density. In this case, a high rate of false positives can occur. (c) The proposed strategy considers learning two hyperplanes with the same normal vector but with different offsets. (d) These hyperplanes learn the “normal” region for the SVM scores. This region is called a slab.

positive rate. Unlike the OCSVM, the proposed OCSSVM approach encloses the normal region of the target class in feature space by using two parallel hyperplanes. When an instance falls inside the normal region or the slab created by the hyperplanes, the OCSSVM labels it as a sample from the target class, and negative otherwise. Figure 2 provides an overview of this new algorithm.

Using two parallel hyperplanes has been explored before in visual recognition problems. Cevikalp and Triggs [4] proposed a cascade of classifiers for object detection. Similarly, Scheirer *et al.* [20] proposed the 1-vs-Set SVM, where a greedy algorithm calculates the slab parameters after training a regular linear SVM. However, these methods are not strictly one-class classifiers since they use samples from known negative classes. Parallel hyperplanes have also been used by Giesen *et al.* [11] to compress a set of 3D points and by Glazer *et al.* [12] to estimate level sets from a high-dimensional distribution. In contrast to these methods, the OCSSVM targets the open-set recognition problem directly and computes the optimal size of the slab automatically.

This work presents two experiments on two publicly available visual recognition datasets. This is because visual recognition systems encounter novel classes very frequently in natural scenes that contain both target and novel objects. The experiments evaluate the performance of the proposed approach and compare it with other state-of-the-art one-class classifiers. The experiments show that OCSSVM consistently outperforms the one-class SVM and performs comparable to or better than other one-class classifiers.

The OCSSVM represents a step towards the ideal robust recognition system based on an ensemble of one-class classifiers. The proposed OCSSVM can also improve the performance of other applications such as the identification of abnormal episodes in gas-turbines [7]; the detection of abnormal medical states from vital signs [14]; and the detection of impostor patterns in a biometric system [14].

#### A. Brief Review of the One-class SVM

Schölkopf *et al.* [22] proposed the one-class support vector machine (OCSVM) to detect novel or outlier samples. Their goal was to find a function that returns +1 in a “small” region capturing most of the target data points, and -1 elsewhere.

Their strategy consists of mapping the data to a feature space via kernel methods. Subsequently, it finds a hyperplane in this new feature space that maximizes the margin between the origin and the data.

To find this hyperplane, Schölkopf *et al.* proposed the following optimization problem:

$$\begin{aligned} & \underset{\mathbf{w}, \rho, \xi}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu m} \sum_{i=1}^m \xi_i - \rho \\ & \text{subject to} && \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle \geq \rho - \xi_i, \\ & && \xi_i \geq 0, \quad i = 1, \dots, m, \end{aligned} \quad (1)$$

where  $m$  is the number of total training samples from the target class;  $\nu$  is an upper-bound on the fraction of outliers and a lower bound on the fraction of support vectors (SV);  $\mathbf{x}_i$  is the  $i$ -th training sample feature vector;  $\mathbf{w}$  is the hyperplane normal vector;  $\Phi()$  is a feature map;  $\xi$  are slack variables; and  $\rho$  is the offset (or threshold). Solvers for this problem compute the dot product  $\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$  via a kernel function  $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$ .

Schölkopf *et al.* [22] proposed to solve the problem shown in Eq. (1) via its dual problem:

$$\begin{aligned} & \underset{\alpha}{\text{minimize}} && \frac{1}{2} \alpha^T K \alpha \\ & \text{subject to} && \|\alpha\|_1 = 1, \\ & && 0 \leq \alpha_i \leq \frac{1}{\nu m}, \quad i = 1, \dots, m, \end{aligned} \quad (2)$$

where  $K$  is the kernel matrix calculated using a kernel function, *i.e.*,  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ , and  $\alpha$  are the dual variables. This optimization problem is a constrained quadratic-program which is convex. Thus, solvers can use Newton-like methods [3], [21] or a variant of the sequential-minimal-optimization (SMO) technique [16].

The SVM decision function is calculated as follows:

$$f(\mathbf{x}) = \text{sgn} \left( \underbrace{\sum_{i=1}^m \alpha_i k(\mathbf{x}_i, \mathbf{x})}_{\langle \mathbf{w}, \Phi(\mathbf{x}) \rangle} - \rho \right), \quad (3)$$

where the offset  $\rho$  can be recovered from the support vectors that lie exactly on the hyperplane, *i.e.*, the training feature

vectors whose dual variables satisfy  $0 < \alpha_i < \frac{1}{\nu m}$ . In this work, the projection  $\langle \mathbf{w}, \Phi(\mathbf{x}) \rangle$  of a sample  $\mathbf{x}$  onto the normal vector  $\mathbf{w}$  is called the SVM score.

### B. Discussion

An interpretation of the solution  $(\mathbf{w}^*, \rho^*)$  for the problem stated in Eq. (1) is a hyperplane that bounds the SVM scores from below; see the inequality constraints in Eq. (1). This interpretation also considers that the SVM score is a random variable. In this context,  $\rho^*$  is a threshold that discards outliers falling on the left tail of the SVM score density. Figures 2(a) and 2(b) illustrate this rationale.

However, the one-class SVM does not account for outliers that occur on the right tail of the SVM-score density. It needs to account for them to reduce false positives. Its decision rule considers these outliers as target samples yielding undesired false positives and decrease of performance.

The proposed strategy does account for these outliers. It learns two hyperplanes that tightly enclose the normal support of the SVM score density from the positive class. These hyperplanes bound the density from “below” and from “above.” The proposed strategy considers samples falling in between these hyperplanes the “normal” state of the positive class SVM scores. It considers samples falling outside these hyperplanes outliers: novel or abnormal samples. The region in between the hyperplanes is called a “slab.” In contrast with the SVM’s default strategy, the proposed strategy assumes that samples from the negative class can have both negative and positive SVM scores; Figures 2(c) and 2(d) illustrate the proposed strategy.

## II. ONE-CLASS SLAB SUPPORT VECTOR MACHINE

This section describes the proposed one-class slab support vector machine. OCSSVM requires two hyperplanes to classify instances as negative (novel or abnormal samples) or positive (target class samples). Both hyperplanes are characterized by the same normal vector  $\mathbf{w}$ , and two offsets  $\rho_1$  and  $\rho_2$ .

The goal of OCSSVM is to find two hyperplanes that tightly enclose the region in feature space of the SVM-score density for the positive class. The positive side of each hyperplane coincides with the slab region and their negative side indicates the area where novel or abnormal samples occur; Figs. 2(c) and 2(d) illustrate the proposed configuration of the hyperplanes and decision process.

OCSSVM solves a convex optimization problem to find the hyperplane parameters  $(\mathbf{w}, \rho_1, \rho_2)$ . This problem is stated as follows:

$$\begin{aligned} & \underset{\mathbf{w}, \rho_1, \rho_2, \xi, \bar{\xi}}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu_1 m} \sum_{i=1}^m \xi_i - \rho_1 + \frac{\varepsilon}{\nu_2 m} \sum_{i=1}^m \bar{\xi}_i + \varepsilon \rho_2 \\ & \text{subject to} && \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle \geq \rho_1 - \xi_i, \xi_i \geq 0, \\ & && \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle \leq \rho_2 + \bar{\xi}_i, \bar{\xi}_i \geq 0, \quad i = 1, \dots, m, \end{aligned} \quad (4)$$

where  $(\mathbf{w}, \rho_1)$  are the parameters for the “lower” hyperplane  $f_1$ ;  $(\mathbf{w}, \rho_2)$  are the parameters of the “upper” hyperplane  $f_2$ ,  $\xi$  and  $\bar{\xi}$  are slack variables for the lower and upper hyperplanes,

respectively;  $\Phi()$  is the implicit feature map in the kernel function; and  $\nu_1$ ,  $\nu_2$ , and  $\varepsilon$  are parameters. The parameter  $\varepsilon$  controls the contribution of the slack variables  $\bar{\xi}$  and the offset  $\rho_2$  to the objective function. The parameters  $\nu_1$  and  $\nu_2$  control the size of the slab.

This proposed optimization problem extends the formulation introduced by Schölkopf *et al.* [22]. It adds two new linear inequality constraints per training sample, which are the constraints for the hyperplane  $f_2$ , and penalty terms in the objective function of the optimization problem shown in Eq. (1). This extension is mainly composed of linear terms and constraints. Consequently, it preserves convexity.

The offsets  $\rho_1$  and  $\rho_2$  have the following interpretation: they are thresholds that bound the SVM scores from the positive class (*i.e.*,  $\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle$ ) from below and above, respectively. This new interpretation motivates the names for the lower and upper hyperplanes mentioned earlier. The region in between these bounds is the “slab,” and its size can be controlled by  $\nu_1$  and  $\nu_2$ . The slack variables  $\xi$  and  $\bar{\xi}$  allow the OCSSVM to exclude some SVM scores that deviate from the slab region: the normal region of the SVM score density from the positive class.

The decision function of the OCSSVM,

$$f(\mathbf{x}) = \text{sgn} \{ (\langle \mathbf{w}, \Phi(\mathbf{x}) \rangle - \rho_1) (\rho_2 - \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle) \}, \quad (5)$$

is positive when SVM scores fall inside the slab region, and negative otherwise.

Solving the primal problem (shown in Eq. (4)) is challenging – especially when a non-linear kernel function is used. However, the dual problem of several SVMs often yields a simpler-to-solve optimization problem. The dual problem for the OCSSVM is

$$\begin{aligned} & \underset{\alpha, \bar{\alpha}}{\text{minimize}} && \frac{1}{2} (\alpha - \bar{\alpha})^T K (\alpha - \bar{\alpha}) \\ & \text{subject to} && 0 \leq \alpha_i \leq \frac{1}{\nu_1 m}, \sum_{i=1}^m \alpha_i = 1, \\ & && 0 \leq \bar{\alpha}_i \leq \frac{\varepsilon}{\nu_2 m}, \sum_{i=1}^m \bar{\alpha}_i = \varepsilon, \quad i = 1, \dots, m, \end{aligned} \quad (6)$$

where  $K$  is the kernel matrix;  $\alpha_i$  and  $\bar{\alpha}_i$  are the  $i$ -th entries for the dual vectors  $\alpha$  and  $\bar{\alpha}$ , respectively; and  $0 \leq \nu_1 \leq 1$ ,  $0 \leq \nu_2 \leq 1$ , and  $0 \leq \varepsilon$  are parameters. This dual problem is a constrained quadratic program that can be solved with convex solvers. This work considers only positive definite kernels, *i.e.*,  $K$  is positive definite [21]. Therefore,  $\varepsilon \neq 1$  must hold to avoid the trivial solution:  $\alpha = \bar{\alpha}$ .

The decision function can be re-written in terms of only the dual variables  $\alpha$ ,  $\bar{\alpha}$  as follows:

$$f(\mathbf{x}) = \text{sgn} \{ (s_{\mathbf{w}} - \rho_1) (\rho_2 - s_{\mathbf{w}}) \}, \quad (7)$$

where

$$s_{\mathbf{w}} = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle = \sum_{i=1}^m (\alpha_i - \bar{\alpha}_i) k(\mathbf{x}, \mathbf{x}_i); \quad (8)$$

and

$$\rho_1 = \frac{1}{N_{SV_1}} \sum_{i: 0 < \alpha_i < \frac{1}{\nu_1 m}}^{N_{SV_1}} \sum_j^m (\alpha_j - \bar{\alpha}_j) k(\mathbf{x}_i, \mathbf{x}_j) \quad (9)$$

$$\rho_2 = \frac{1}{N_{SV_2}} \sum_{i: 0 < \bar{\alpha}_i < \frac{\varepsilon}{\nu_2 m}}^{N_{SV_2}} \sum_j^m (\alpha_j - \bar{\alpha}_j) k(\mathbf{x}_i, \mathbf{x}_j). \quad (10)$$

The SVM score  $s_{\mathbf{w}}$  is obtained from Eq. (8) and re-writing dot products with the kernel function. On the other hand, the offsets require analysis from the KKT conditions (see Appendix A) to establish their relationship with the dual variables. The offset computation requires knowledge of the support vectors that lie exactly on the lower and upper hyperplanes. These support vectors are detected by evaluating if their dual variables satisfy  $0 < \alpha_i < \frac{1}{\nu_1 m}$  and  $0 < \bar{\alpha}_i < \frac{\varepsilon}{\nu_2 m}$  for the lower and upper hyperplane, respectively. Equations (9) and (10) require the number of support vectors  $N_{SV_1}$ ,  $N_{SV_2}$  that exactly lie on the lower and upper hyperplanes, respectively. Moreover, it can be shown via the KKT conditions that if  $\alpha_i > 0$ , then  $\bar{\alpha}_i = 0$ , and that if  $\bar{\alpha}_i > 0$ , then  $\alpha_i = 0$ . This means that each hyperplane has its own set of support vectors; the reader is referred to the Appendix A for a more detailed analysis of the KKT conditions.

### III. EXPERIMENTS

This section presents two experiments (described in Sections III-A and III-B) that assess the performance of the proposed OCSSVM. These experiments use two different publicly available datasets: the letter dataset [10] and the PascalVOC 2012 [9] dataset.

We implemented a primal-dual interior point method solver in C++<sup>1</sup> to find the hyperplane parameters of the proposed OCSSVM. The experiments on the letter dataset were carried out on a MacBook Pro with 16GB of RAM and an Intel core i7 CPU. The experiments on the PascalVOC 2012 dataset were executed on a machine with 32GB of RAM and an Intel core i7 CPU.

The experiments compared the proposed approach to other state-of-the-art one-class classifiers: support vector data description (SVDD) [23], one-class kernel PCA (KPCA) [13], kernel density estimation (KDE), and the one-class support vector machine (OCSVM) [22] – the main baseline. The experiments used the implementations from LibSVM [5] for SVDD and SVM; and a publicly available Matlab implementation we created for the one-class kernel PCA algorithm to apply to the letter dataset. However, the experiments used a C++ KPCA implementation (also developed in house) for the PascalVOC 2012 dataset, since the Matlab implementation struggled with the high dimensionality of the feature vectors and large number of samples in the dataset. For the multivariate kernel density estimation, we used Ihler’s publicly available Matlab toolkit<sup>2</sup>. However, the KDE method did not run on the PascalVOC 2012 dataset due to the large volume of data. Thus, the experiments omit KDE results for that dataset.

<sup>1</sup><http://vfragoso.com>

<sup>2</sup>Multivariate KDE: <http://www.ics.uci.edu/~ihler/code/kde.html>

TABLE I  
MEDIAN MATTHEWS CORRELATION COEFFICIENTS OVER THE 26 LETTERS. BOLD NUMBERS INDICATE THE HIGHEST SCORE FOR A KERNEL (ROW). OCSSVM CONSISTENTLY OUTPERFORMED OCSVM AND PERFORMED COMPARABLE TO OR BETTER THAN THE REMAINING ONE-CLASS CLASSIFIERS.

Kernel	KDE	KPCA	SVDD	OCSVM	OCSSVM
Linear	-	0.01	0.09	0.02	<b>0.14</b>
RBF	0.18	0.17	0.11	0.07	<b>0.39</b>
Intersection	-	0.18	0.01	0.04	<b>0.26</b>
Hellinger	-	0.01	0.02	0.02	<b>0.13</b>
$\chi^2$	-	<b>0.18</b>	0.02	0.02	<b>0.18</b>

The experiments trained a one-class classifier for each class in the datasets. Recall that one-class classifiers only use positive samples for training. To evaluate the performance of the one-class classifiers, the experiments used the remaining classes as negative samples (*i.e.*, novel class instances). The tested datasets are unbalanced in this setting since there are more instances from the negative class compared to the positive class. Note that common metrics such as precision, recall, and f1-measure are sensitive to unbalanced datasets. This is because they depend on the counts of true positives, false positives, and false negatives.

Fortunately, the Matthews correlation coefficient (MCC) [17] is known to be robust to unbalanced datasets. The MCC ranges between  $-1$  and  $+1$ . A coefficient of  $+1$  corresponds to perfect prediction,  $0$  corresponds to an equivalent performance of random classification, and  $-1$  corresponds to a perfect disagreement between predictions and ground truth labels; see Appendix D material for more details about MCC.

The experiment used common kernels (*e.g.*, linear and radial basis function (RBF)) as well as efficient additive kernels [24] (*e.g.*, intersection, Hellinger, and  $\chi^2$ ). Among these kernels, only the RBF kernel requires setting a free parameter:  $\gamma$ . Also, the experiment used a Gaussian kernel for the KDE method. Its bandwidth was determined by the rule-of-thumb method, an automatic algorithm for kernel bandwidth estimation included in the used Matlab KDE toolbox. The experiments compare the KDE method only with the remaining one-class classifiers using an RBF kernel since the Gaussian kernel belongs to that family.

The experiments ran a grid-search over various kernel and classifier parameters, such as  $\gamma$  for the RBF kernel,  $C$  parameter for SVDD,  $\nu_1, \nu_2, \nu$  for the one-class SVMs, and number of components for KPCA, using a validation set for every class in every dataset; the reader is referred to the Appendix C where these parameters are shown.

To determine the  $\varepsilon$  parameters for training the proposed OCSSVM, the experiments used a toy dataset where samples from a bivariate Normal distribution were used. It was observed that  $\varepsilon = \frac{2}{3}$  produced good results; see Appendix B for more details of this process.

#### A. Evaluation on Letter Dataset

This experiment aims at evaluating the performance of the OCSSVM. The tested dataset is letter [10], which contains 20,000 feature vectors of the 26 capital letters in the English

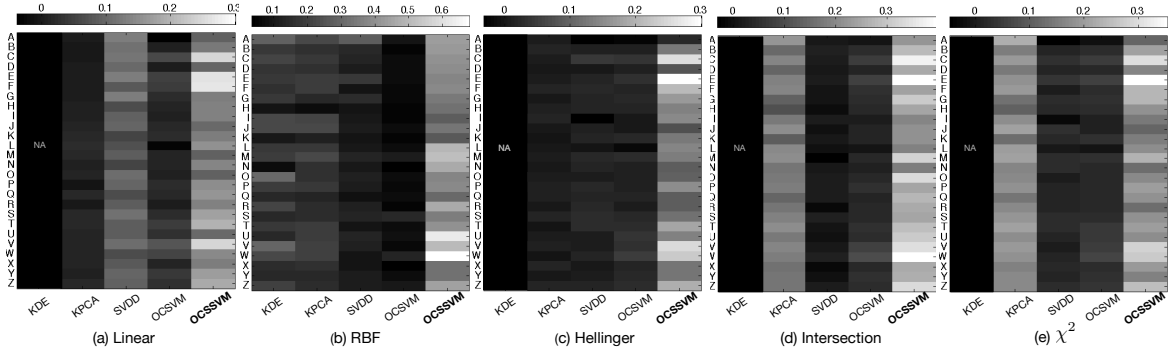


Fig. 3. Matthews correlation coefficient on the letter dataset across different kernels (a-e); brighter indicates better performance. The proposed OCSSVM performed comparable or better than one-class kernel PCA (KPCA), kernel density estimation (KDE), support vector data description (SVDD), and one-class SVM (OCSVM). A comparison with the KDE method is only valid when using the RBF kernel.

alphabet. Each feature vector is a 16-dimensional vector capturing statistics of a single character. The dataset provides 16,000 samples for training and 4,000 for testing. The one-class classification problem consists of training the classifier with instances of a single character (the positive class), and detecting instances of that character in the presence of novel classes – instances of the remaining 25 characters.

Figure 3 shows the results of this experiment. It visualizes the performance of the tested classifiers across classes for different kernels. Table I presents a performance summary per kernel and per method. The results shown in Figure 3 and Table I only include a comparison of the KDE method and the one-class classifiers with an RBF kernel since the KDE method uses a Gaussian kernel, which belongs to the RBF family. Because the experiment uses Matthews correlation coefficient (MCC), higher scores imply better performance. Thus, a consistent bright vertical stripe in a visualization indicates good performance across all the classes in the dataset for a particular kernel. The figure shows that the proposed OCSSVM tends to have a consistent bright vertical stripe across different kernels and classes. This can be confirmed in Table I where OCSSVM achieves the highest median MCC for all of the kernels. The visualizations also show that the proposed OCSSVM outperformed the SVM method consistently. Comparing the OCSSVM and the SVM columns in Table I confirms the better performance of the proposed method. Table I also shows that OCSSVM performed comparable or better than one-class kernel PCA (KPCA), kernel density estimation (KDE), and support vector data description (SVDD).

#### B. Evaluation on PascalVOC 2012 Dataset

The goal of this experiment is to assess the performance of the OCSSVM on a more complex dataset: PascalVOC 2012 [9]. This dataset contains 20 different visual classes (objects) and provides about 1,000 samples per class. It has been used mainly for object detection. The experiment used HOG [8] features for every object class. To mimic novel classes that an object detector encounters, the experiment randomly picked 10,000 background regions for which HOG features were computed. The dimensionality of these features per class ranges from 2,304 to 36,864. This experiment used high-dimensional feature vectors and a large number of samples. Consequently, the kernel density estimation (KDE)

TABLE II  
MEDIAN OF THE 3-FOLD MATTHEWS CORRELATION COEFFICIENTS OVER THE 20 CLASSES IN THE PASCALVOC 2012 DATASET PER KERNEL. BOLD NUMBERS INDICATE THE HIGHEST SCORE FOR A KERNEL (ROW). OCSSVM OUTPERFORMED THE OCSVM IN MOST OF THE CASES, WITH THE EXCEPTION OF THE RBF KERNEL CASE. IT PERFORMED COMPARABLE TO OR BETTER THAN ONE-CLASS KPCA AND SVDD.

Kernel	KPCA	SVDD	OCSVM	OCSSVM
Linear	0.02	<b>0.09</b>	0.01	0.07
RBF	0.05	0.07	<b>0.14</b>	0.09
Intersection	0.18	0.01	0.04	<b>0.26</b>
Hellinger	0.01	0.02	0.02	<b>0.13</b>
$\chi^2$	<b>0.18</b>	0.02	0.02	<b>0.18</b>

MATLAB toolkit struggled and did not run properly on this dataset. Hence, the experiment omits the result for this method.

The experiment trained one-class classifiers for each object using a 3-fold cross-validation procedure. The testing set for a fold was composed of object samples and all background features. Figure 4 shows the visualizations of the average Matthews correlation coefficients (MCC) for this experiment. In addition, Table II presents a summary of this experiment.

Table II shows that OCSSVM tended to outperform the one-class SVM across kernels. Moreover, it performed comparable to or better than one-class KPCA and SVDD across kernels. Figure 4 shows that the OCSSVM tended to outperform the SVM method across classes and kernels.

#### IV. CONCLUSIONS AND FUTURE DIRECTIONS

This work presented the one-class slab support vector machine as a step towards the idealized one-class solution for open-set recognition. In contrast to the regular one-class SVM, which learns a single hyperplane for identifying target samples, instances from the positive class, the proposed classifier uses two parallel hyperplanes learned in feature space to enclose a portion of the target samples. However, each plane has an offset with respect to the origin that places them at different locations in feature space, creating a “slab.” The proposed approach to train the OCSSVM is a quadratic program (QP) that estimates the hyperplane normal vector and the two offsets.

The proposed OCSSVM showed consistent performance improvement over the regular one-class SVM on two different datasets: letter [10] and the PascalVOC 2012 [9]. The proposed strategy performed comparable or better than other state of

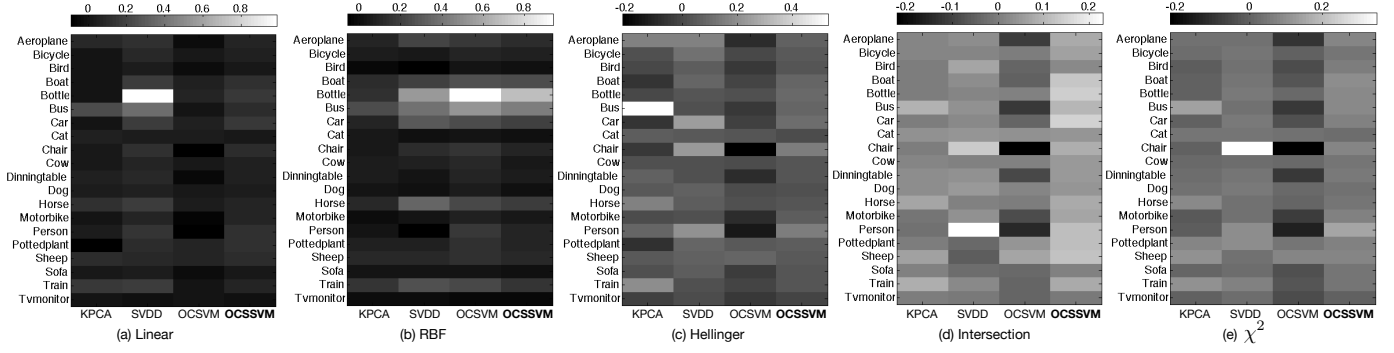


Fig. 4. Average of the 3-fold Matthews correlation coefficient scores per class; brighter indicates better performance. The proposed OCSSVM outperformed the SVM using efficient additive kernels (Hellinger, Intersection, and  $\chi^2$ ). It performed comparable or better than the one-class kernel PCA (KPCA), the support vector data description (SVDD), and the one-class SVM (OCSVM).

the art one-class classifiers, such as support vector data description [23], one-class kernel PCA [13], and kernel density estimation.

The approach used a Newton-based QP solver to train the OCSSVM. However, this solver is not efficient and a derivation of a sequential-minimal-optimization (SMO) [16] is planned for future work. The plan includes the adaptation of the SMO solver to deal with an extra inequality constraint that the QP of the OCSSVM includes.

#### APPENDIX A KKT ANALYSIS

In this section we explore the different cases that the optimal values for the dual variables  $\alpha$  and  $\bar{\alpha}$  can fall in. As a result of this analysis, we learn how to obtain the offset values  $\rho_1$  and  $\rho_2$ , useful conditions on support vectors for each hyperplane, and invalid cases. To do so, we exploit the KKT conditions at their optimal values. The optimal dual variables must satisfy the following statements:

$$\begin{cases} \alpha_i (\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle - \rho_1 + \xi_i) = 0 \\ \beta_i \xi_i = 0 \\ \bar{\alpha}_i (\rho_2 + \bar{\xi}_i - \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle) = 0 \\ \bar{\beta}_i \bar{\xi}_i = 0 \end{cases}.$$

Before starting to analyze the cases, we need to remember the following relationships:

$$\mathbf{w} = \sum_{i=0}^m (\alpha_i - \bar{\alpha}_i) \Phi(\mathbf{x}) \quad (11)$$

$$\beta_i = \frac{1}{\nu_1 m} - \alpha_i \quad (12)$$

$$1 = \sum_i \alpha_i \quad (13)$$

$$\bar{\beta}_i = \frac{\varepsilon}{\nu_2 m} - \bar{\alpha}_i \quad (14)$$

$$\varepsilon = \sum_i \bar{\alpha}_i, \quad (15)$$

which are obtained by differentiating the Laplacian of our problem shown in Eq. (6) of the main submission.

#### A. Cases

- 1) Case  $\alpha_i = 0$  and  $\bar{\alpha}_i = 0$ . Given this scenario we conclude using Equations (12), (13), (14), and (15) that

$$\begin{aligned} \beta_i &= \frac{1}{\nu_1 m} \\ \bar{\beta}_i &= \frac{\varepsilon}{\nu_2 m} \end{aligned} \quad (16)$$

Therefore,

$$\xi_i = \bar{\xi}_i = 0. \quad (17)$$

This implies that there are no slack variables compensating for the inequalities in the primal problem shown in Eq. (4) and thus we conclude that

$$\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle > \rho_1 \quad (18)$$

$$\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle < \rho_2. \quad (19)$$

Samples with  $\alpha_i = 0$  and  $\bar{\alpha}_i = 0$  are instances that fall inside the slab.

- 2) Case  $0 < \alpha_i < \frac{1}{\nu_1 m}$  and  $\bar{\alpha}_i = 0$ . In this case

$$\begin{cases} \beta_i = \frac{1}{\nu_1 m} - \alpha_i > 0 \\ \bar{\beta}_i = \frac{\varepsilon}{\nu_2 m} \end{cases} \quad (20)$$

Therefore, the following must be true

$$\begin{cases} \xi_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle = \rho_1 \\ \bar{\xi}_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle < \rho_2. \end{cases} \quad (21)$$

- 3) Case  $\alpha_i = 0$  and  $0 < \bar{\alpha}_i < \frac{1}{\nu_1 m}$ . In this case

$$\begin{cases} \bar{\beta}_i = \frac{\varepsilon}{\nu_2 m} - \bar{\alpha}_i > 0 \\ \beta_i = \frac{1}{\nu_1 m} \end{cases} \quad (22)$$

Therefore, the following must be true

$$\begin{cases} \xi_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle = \rho_2 \\ \bar{\xi}_i = 0. \end{cases} \quad (23)$$

- 4) Case  $0 < \bar{\alpha}_i < \frac{1}{\nu_1 m}$  and  $0 < \alpha_i < \frac{\varepsilon}{\nu_1 m}$ . This implies that

$$\begin{cases} \bar{\beta}_i = \frac{\varepsilon}{\nu_2 m} - \bar{\alpha}_i > 0 \\ \beta_i = \frac{1}{\nu_1 m} - \alpha_i > 0 \end{cases} \quad (24)$$

Therefore,

$$\begin{cases} \xi_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle = \rho_2 \\ \bar{\xi}_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle = \rho_1 \end{cases} \quad (25)$$

Note that this case by construction of the primal problem should not happen. This case implies that the size of the slab (*i.e.*,  $\rho_2 - \rho_1$ ) is zero. In other words, the two planes overlap. Therefore, there is no slab in the feature space and by construction this should not happen.

- 5) Case  $\alpha_i = \frac{1}{\nu_1 m}$  and  $\bar{\alpha}_i = 0$ . This situation implies that

$$\begin{cases} \bar{\beta}_i = \frac{\varepsilon}{\nu_2 m} \\ \beta_i = 0 \end{cases} \quad (26)$$

Therefore, we conclude that

$$\begin{cases} \xi_i > 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle < \rho_2 \\ \bar{\xi}_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle < \rho_1 \end{cases} \quad (27)$$

Another implication of this case is that the  $i$ -th sample is considered an outlier/novel sample with respect to the first plane.

- 6) Case  $\bar{\alpha}_i = \frac{\varepsilon}{\nu_2 m}$  and  $\alpha_i = 0$ . This case implies that

$$\begin{cases} \bar{\beta}_i = 0 \\ \beta_i = \frac{1}{\nu_1 m} \end{cases} \quad (28)$$

Therefore, we conclude that

$$\begin{cases} \xi_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle > \rho_2 \\ \bar{\xi}_i > 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle > \rho_1 \end{cases} \quad (29)$$

Again, the  $i$ -th sample is considered an outlier/novel sample with respect to the second plane.

- 7) Case  $\bar{\alpha}_i = \frac{\varepsilon}{\nu_2 m}$  and  $0 < \alpha_i < \frac{1}{\nu_1 m}$ . In this case we have

$$\begin{cases} \bar{\beta}_i = 0 \\ \beta_i = \frac{1}{\nu_1 m} - \alpha_i > 0 \end{cases} \quad (30)$$

Therefore,

$$\begin{cases} \xi_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle > \rho_2 \\ \bar{\xi}_i > 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle = \rho_1 \end{cases} \quad (31)$$

This implies that  $\rho_2 < \rho_1$ , which again, by construction cannot happen. Thus, this case must not occur.

- 8) Case  $\alpha_i = \frac{1}{\nu_1 m}$  and  $0 < \bar{\alpha}_i < \frac{\varepsilon}{\nu_2 m}$ . In this case we have

$$\begin{cases} \beta_i = 0 \\ \bar{\beta}_i = \frac{\varepsilon}{\nu_2 m} - \bar{\alpha}_i > 0 \end{cases} \quad (32)$$

Therefore,

$$\begin{cases} \xi_i > 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle = \rho_2 \\ \bar{\xi}_i = 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle < \rho_1 \end{cases} \quad (33)$$

This implies that  $\rho_2 < \rho_1$ , which again, by construction cannot happen. Thus, this case must not occur.

- 9) Case  $\bar{\alpha}_i = \frac{\varepsilon}{\nu_2 m}$  and  $\alpha_i = \frac{1}{\nu_1 m}$ . This implies that

$$\begin{cases} \beta_i = \frac{1}{\nu_1 m} - \alpha_i > 0 \\ \bar{\beta}_i = \frac{\varepsilon}{\nu_2 m} - \bar{\alpha}_i > 0 \end{cases} \quad (34)$$

Therefore,

$$\begin{cases} \xi_i > 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle > \rho_2 \\ \bar{\xi}_i > 0 \\ \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle < \rho_1 \end{cases} \quad (35)$$

This scenario implies that  $\rho_2 < \rho_1$ , which again, contradicts our construction of the problem. Therefore this must not occur.

We can conclude from the analysis of these cases that any plane contains the  $i$ -th sample when its corresponding dual satisfies  $0 < \alpha_i < \frac{1}{\nu_1 m}$  or  $0 < \bar{\alpha}_i < \frac{\varepsilon}{\nu_2 m}$  for the lower and higher hyperplanes, respectively. However, only one plane can contain the  $i$ -th sample at a time. Therefore, at the optimal point  $\alpha_i > 0$  and  $\bar{\alpha}_i > 0$  does not occur. It only happens exclusively.

Thus, to recover the offsets  $\rho_1$  and  $\rho_2$  we need to collect all the points that satisfy either  $0 < \alpha_i < \frac{1}{\nu_1 m}$  or  $0 < \bar{\alpha}_i < \frac{\varepsilon}{\nu_2 m}$ . Thus,

$$\rho_1 = \frac{1}{n_1} \sum_{i: 0 < \alpha_i < \frac{1}{\nu_1 m}} \langle \mathbf{w}, \Phi(x_i) \rangle, \quad (36)$$

where  $n_1$  is the number of points that satisfy  $0 < \alpha_i < \frac{1}{\nu_1 m}$ . In a similar fashion, we can recover offset  $\rho_2$ :

$$\rho_2 = \frac{1}{n_2} \sum_{i: 0 < \bar{\alpha}_i < \frac{\varepsilon}{\nu_2 m}} \langle \mathbf{w}, \Phi(x_i) \rangle, \quad (37)$$

where  $n_2$  is the number of points that satisfy  $0 < \bar{\alpha}_i < \frac{\varepsilon}{\nu_2 m}$ .

## APPENDIX B TOY DATASET EXPERIMENTS

The goal of this experiment is twofold: 1) obtain insight about our proposed method and visualize the computed decision function for two kernels: linear and radial basis function (RBF); and 2) explore the effect of  $\varepsilon$  on the learned hyperplanes.



TABLE III

FRACTION OF POINTS THAT THE ONE-CLASS SLAB SVM CONSIDERS AS POSITIVE SAMPLES AS A FUNCTION OF  $\varepsilon$ . THE FRACTION OF POINTS LABELED AS POSITIVE SAMPLES DID NOT CHANGE SIGNIFICANTLY REGARDLESS OF THE KERNEL AND THE VALUE OF  $\varepsilon$ .

Kernel	$\varepsilon = 1/6$	$\varepsilon = 2/6$	$\varepsilon = 3/6$	$\varepsilon = 4/6$	$\varepsilon = 5/6$
Linear	0.91	0.91	0.91	0.91	0.92
RBF	0.90	0.90	0.90	0.90	0.90

### A. Parameter Exploration

The goal of this experiment is to determine a good value for the  $\varepsilon$  parameter. To do so we generated a toy dataset composed of 1500 points drawn from a bivariate Normal distribution. We trained our one-class slab SVM using a linear kernel and an RBF kernel with  $\gamma = 0.5$ , with  $\nu_1 = 0.1$  and  $\nu_2 = 0.05$ .

We varied the values of  $\varepsilon$  in the interval  $[\frac{1}{6}, \frac{5}{6}]$ . A visualization of the hyperplanes is shown in Fig. 5 and Fig. 6. The visualizations show that there is no significant differences in the learned hyperplanes when  $\varepsilon$  is varied across kernels. To verify this, we calculated the fraction of points that were considered positive by each of the learned hyperplanes. The results are shown in Table III. Thus we conclude that the value of  $\varepsilon$  does not affect significantly the learned hyperplanes.

### B. Insight About One-Class Slab SVM

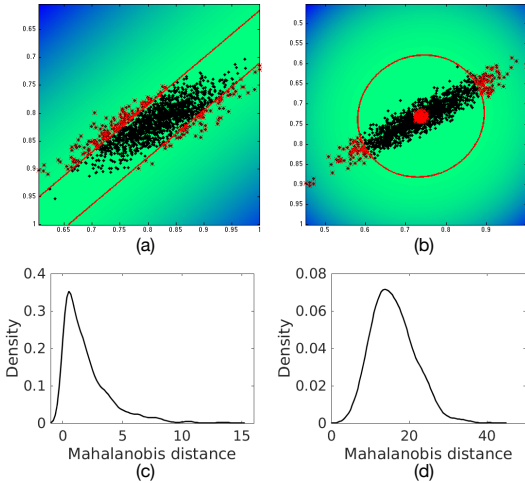


Fig. 7. One-class slab SVM decision functions on a toy dataset. The support vectors as well as the hyperplanes are shown in red. (a) The computed slab using a linear kernel encloses most of the bivariate Normal points. (b) The “doughnut” like slab computed using a radial basis function (RBF) kernel captures two sets of extreme points: points deviating from the norm, and points very close to the mean. (c) The extremes found by the RBF kernel can be explained via the density of the Mahalanobis distance between the mean and a point in the dataset. It is very unlikely to observe a point very close to the mean. (d) The chances of observing a point close to the mean becomes less unlikely when the dimensionality of the points increases. This can be seen by observing the Mahalanobis distance between the mean and a point in the dataset with dimensionality 16.

For this experiment we set  $\varepsilon = \frac{2}{3}$ ,  $\nu_1 = 0.1$ , and  $\nu_2 = 0.05$ . Our toy dataset is composed of 1500 points drawn from a bivariate Normal distribution. We trained our one-class slab SVM using a linear and an RBF kernel with  $\gamma = 0.5$ . We show a visualization of the computed decision functions in Fig. 7. The linear kernel finds a slab in the input space that

TABLE V

RBF KERNEL PARAMETER ( $\gamma$ ) FOR THE PASCALVOC DATASET.

Aeroplane	Bicycle	Bird	Boat	Bottle
9.5367e-07	9.5367e-07	9.5367e-07	9.5367e-07	3.8147e-06
Bus	Car	Cat	Chair	Cow
2.3842e-07	9.5367e-07	2.3842e-07	1.1921e-07	9.5367e-07
Diningtable	Dog	Horse	Motorbike	Person
4.7684e-07	1.1921e-07	9.5367e-07	4.7684e-07	1.1921e-07
Pottedplant	Sheep	Sofa	Train	Tvmonitor
1.9073e-06	1.9073e-06	2.3842e-07	9.5367e-07	4.7684e-07

captures most of the training data. The RBF kernel finds a slab in the input space that resembles a “doughnut” like slab. The RBF kernel identifies two sets of points that corresponds to the following extremes: 1) points that deviate significantly from the norm; and 2) points that fall very close to the norm. These sets of points can be verified to be “extreme” by analysing the density of the Mahalanobis distance between the mean and a point in the dataset. In Fig. 7(c), not only can we observe that points falling far from the mean are rare, but also points falling very close to the mean are; the peak of the density is close to zero, but it is not exactly zero. This becomes more evident when the dimensionality of points drawn from a multivariate Normal distribution increases; see Fig. 7(d) for an illustration.

## APPENDIX C

### PARAMETERS

In this section we present the parameters we used for the experiments presented in Section 3 of the main submission. These parameters were obtained after running a 5-fold cross validation using a validation set. The criterion was to maximize the recall rate.

#### A. One-class SVM parameters

The  $\nu$  parameter converged to  $\nu = 0.1$  for both datasets. The single kernel that required a parameter to be set, was the RBF kernel. For this kernel we show the parameters used for the letter and PascalVOC datasets in Table IV and Table V, respectively.

#### B. SVDD parameters

The support vector data description (SVDD) method requires a parameter  $C$  for training. We present the  $C$  parameter we used for both experiments and per kernel in Tables VI and VIII. The RBF kernel parameters used for the letter dataset and PascalVOC dataset are shown in Table VII and Table IX, respectively.

#### C. One-class Kernel PCA

The number of components used in both experiments was 16. The RBF kernel parameters ( $\gamma$ ) that we used for the letter and PascalVOC datasets are shown in Table X and Table XI.

#### D. One-Class Slab SVM

The  $\nu_1$  and  $\nu_2$  parameters converged to  $\nu_1 = 0.10$  and  $\nu_2 = 0.01$  for both datasets. The single kernel that required a parameter to be set was the RBF kernel. We show the RBF parameters used for the letter and PascalVOC datasets in Table XII and Table 3 XIII, respectively.



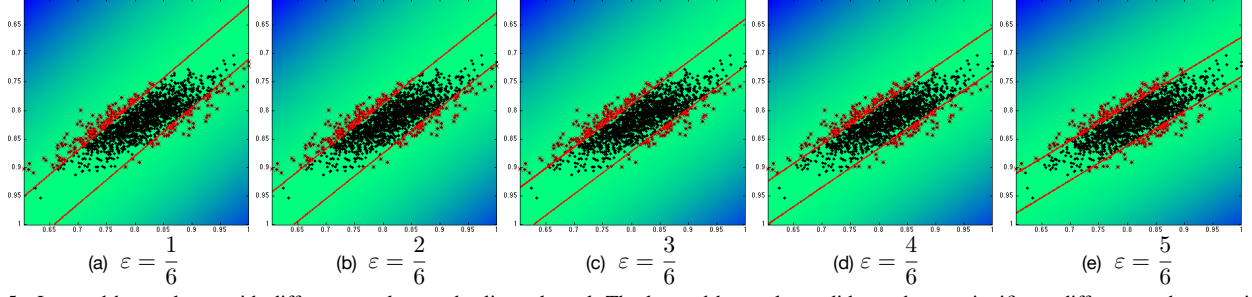


Fig. 5. Learned hyperplanes with different  $\varepsilon$  values and a linear kernel. The learned hyperplanes did not show a significant difference when varying  $\varepsilon$ .

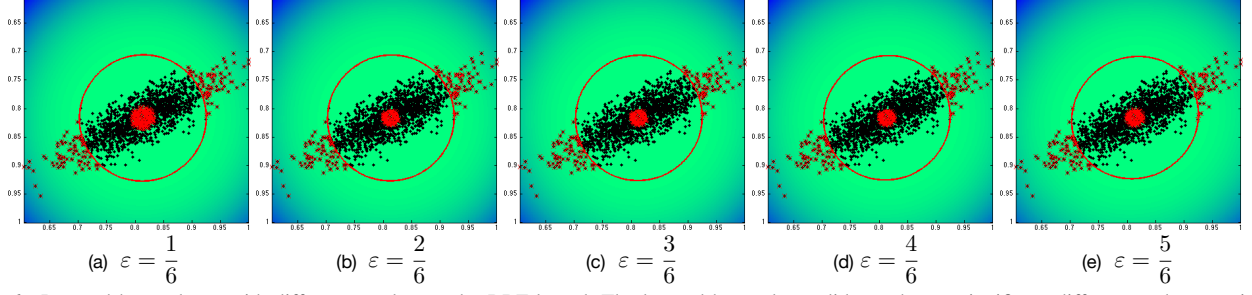


Fig. 6. Learned hyperplanes with different  $\varepsilon$  values and a RBF kernel. The learned hyperplanes did not show a significant difference when varying  $\varepsilon$ .

TABLE IV  
RBF KERNEL PARAMETER ( $\gamma$ ) FOR THE LETTER DATASET.

A	B	C	D	E	F	G	H	I	J	K	L	M
1.0	0.5	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.5
N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1.0	0.5	0.5	2.0	0.5	2.0	1.0	1.0	0.5	1.0	0.5	1.0	1.0

TABLE VI  
SVDD  $C$  PARAMETER FOR THE LETTER DATASET.

Kernel	A	B	C	D	E	F	G	H	I	J	K	L	M
Linear	0.5	0.4	0.4	0.4	0.3	0.3	0.5	0.5	0.4	0.9	0.5	0.5	0.4
RBF	0.9	0.4	0.3	0.2	0.2	0.2	0.7	0.5	0.4	0.2	0.4	0.3	0.8
Intersection	0.1	0.5	0.9	0.3	0.1	0.8	0.6	0.1	0.6	0.3	0.7	0.5	0.5
Hellinger	0.9	0.6	0.9	0.1	0.1	0.1	0.9	0.1	0.8	0.3	0.9	0.2	0.1
$\chi^2$	0.1	0.1	0.6	0.4	0.1	0.1	0.9	0.1	0.6	0.4	0.1	0.1	0.1
Kernel	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
Linear	0.3	0.3	0.5	0.5	0.5	0.5	0.3	0.5	0.4	0.5	0.3	0.4	0.3
RBF	0.5	0.3	0.9	0.3	0.4	0.3	0.4	0.4	0.2	0.5	0.3	0.2	0.6
Intersection	0.5	0.7	0.8	0.2	0.4	0.8	0.8	0.1	0.4	0.1	0.7	0.4	0.3
Hellinger	0.6	0.1	0.1	0.1	0.1	0.4	0.1	0.6	0.1	0.5	0.1	0.1	0.5
$\chi^2$	0.8	0.8	0.9	0.6	0.1	0.9	0.6	0.6	0.1	0.5	0.1	0.1	0.8

TABLE VII  
RBF KERNEL PARAMETER ( $\gamma$ ) FOR SVDD AND THE LETTER DATASET.

A	B	C	D	E	F	G	H	I	J	K	L	M
1.0	0.5	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.5	1.0	0.5
N	O	P	Q	R	S	T	U	V	W	X	Y	Z
0.5	0.5	0.5	1.0	0.5	1.0	1.0	1.0	0.5	1.0	0.5	1.0	0.5

TABLE X  
RBF KERNEL PARAMETER ( $\gamma$ ) FOR ONE-CLASS KERNEL PCA AND THE LETTER DATASET.

A	B	C	D	E	F	G	H	I	J	K	L	M
1.0	0.5	1.0	1.0	4.0	2.0	2.0	16.0	1.0	4.0	16.0	4.0	1.0
N	O	P	Q	R	S	T	U	V	W	X	Y	Z
0.5	0.5	2.0	16.0	4.0	16.0	2.0	8.0	16.0	1.0	2.0	16.0	16.0

TABLE VIII  
SVDD  $C$  PARAMETER FOR THE PASCALVOC DATASET.

Kernel	Aeroplane	Bicycle	Bird	Boat	Bottle
Linear	0.1	0.1	0.1	0.1	0.1
RBF	0.2	0.2	0.2	0.6	0.2
Intersection	0.4	0.1	0.1	0.2	0.7
Hellinger	0.1	0.1	0.7	0.5	0.7
$\chi^2$	0.1	0.9	0.2	0.1	0.9

Kernel	Bus	Car	Cat	Chair	Cow
Linear	0.2	0.1	0.1	0.2	0.1
RBF	0.1	0.2	0.1	0.1	0.2
Intersection	0.2	0.1	0.4	0.1	0.4
Hellinger	0.6	0.1	0.1	0.9	0.1
$\chi^2$	0.2	0.7	0.1	0.1	0.1

Kernel	Dinningtable	Dog	Horse	Motorbike	Person
Linear	0.1	0.2	0.1	0.1	0.1
RBF	0.1	0.2	0.1	0.2	0.1
Intersection	0.4	0.1	0.3	0.1	0.6
Hellinger	0.2	0.1	0.1	0.1	0.7
$\chi^2$	0.1	0.1	0.1	0.1	0.5

Kernel	Pottedplant	Sheep	Sofa	Train	Tvmonitor
Linear	0.1	0.2	0.1	0.1	0.2
RBF	0.1	0.1	0.1	0.1	0.6
Intersection	0.4	0.1	0.2	0.3	0.1
Hellinger	0.1	0.7	0.1	0.1	0.9
$\chi^2$	0.1	0.1	0.1	0.1	0.2

TABLE IX  
RBF KERNEL PARAMETER ( $\gamma$ ) FOR SVDD AND THE PASCALVOC DATASET.

Aeroplane	Bicycle	Bird	Boat	Bottle
4.7684e-07	4.7684e-07	4.7684e-07	7.6294e-06	7.6294e-06
Bus	Car	Cat	Chair	Cow
3.8147e-06	3.8147e-06	1.1921e-07	1.1921e-07	9.5367e-07
Dinningtable	Dog	Horse	Motorbike	Person
2.3842e-07	1.1921e-07	9.5367e-07	1.1921e-07	1.1921e-07
Pottedplant	Sheep	Sofa	Train	Tvmonitor
1.9073e-06	1.9073e-06	1.1921e-07	1.9073e-06	2.3842e-07

#### APPENDIX D MATTHEWS CORRELATION COEFFICIENT

The MCC is computed as follows:

$$MCC = \frac{TP \cdot TN - FN \cdot FP}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}}, \quad (38)$$

where the number of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN) are considered; true and false negatives are the correct and incorrect predictions of negative instances, respectively.

The MCC is positive when the product between  $TN \cdot TP$  is larger than  $FN \cdot FP$ , which only can occur when correct predictions take place. On the other hand, it is negative when the  $FN \cdot FP$  is larger than  $TN \cdot TP$ . The denominator ensures that the MCC metric falls in the  $[-1, +1]$  range. The MCC metric is more robust for unbalanced datasets because the term

TABLE XI  
RBF KERNEL PARAMETER ( $\gamma$ ) FOR ONE-CLASS KERNEL PCA AND THE LETTER DATASET.

Aeroplane	Bicycle	Bird	Boat	Bottle
9.31E-10	2.38E-07	7.45E-09	9.54E-07	9.31E-10
Bus	Car	Cat	Chair	Cow
9.31E-10	2.38E-07	9.31E-10	1.49E-08	7.45E-09
Dinningtable	Dog	Horse	Motorbike	Person
9.31E-10	9.31E-10	9.31E-10	9.31E-10	9.31E-10
Pottedplant	Sheep	Sofa	Train	Tvmonitor
7.63E-06	1.53E-05	9.31E-10	9.31E-10	2.38E-07

TABLE XII  
RBF KERNEL PARAMETER ( $\gamma$ ) FOR THE LETTER DATASET.

A	B	C	D	E	F	G	H	I	J	K	L	M
1.0	0.5	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.5
N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1.0	0.5	0.5	2.0	0.5	2.0	1.0	1.0	0.5	1.0	0.5	1.0	1.0

measuring accurate predictions (*i.e.*,  $TN \cdot TP$ ) considers metrics for both classes. The MCC metric thus measures the overall accuracy of the classifier in a robust manner.

#### REFERENCES

- [1] A. Bendale and T. E. Boulton. Towards open world recognition. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, June 2015.
- [2] C. M. Bishop. Novelty detection and neural network validation. In *IEE Proc. Vision, Image and Signal Processing*, 1994.
- [3] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [4] H. Cevikalp and B. Triggs. Efficient object detection using cascades of nearest convex model classifiers. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2012.
- [5] C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011.
- [6] D. A. Clifton, S. Huguely, and L. Tarassenko. Novelty detection with multivariate extreme value statistics. *Journal of signal processing systems*, 65(3):371–389, 2011.
- [7] D. A. Clifton, L. Tarassenko, N. McGrogan, D. King, S. King, and P. Anuzis. Bayesian extreme value statistics for novelty detection in gas-turbine engines. In *Proc. of the IEEE Aerospace Conf.*, 2008.
- [8] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2005.
- [9] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *Intl. Journal of Computer Vision*, 88(2):303–338, 2010.
- [10] P. W. Frey and D. J. Slate. Letter Recognition Using Holland-Style Adaptive Classifiers. *Machine Learning*, 6(2):161–182, 1991.
- [11] J. Giesen, S. Spalinger, and B. Schölkopf. Kernel methods for implicit surface modeling. In *Advances in Neural Information Processing Systems*, pages 1193–1200, 2004.
- [12] A. Glazer, M. Lindenbaum, and S. Markovitch. q-ocsvm: A q-quantile estimator for high-dimensional distributions. In *Advances in Neural Information Processing Systems*, pages 503–511, 2013.
- [13] H. Hoffmann. Kernel PCA for novelty detection. *Pattern Recognition*, 40(3):863–874, 2007.
- [14] H.-j. Lee and S. Cho. Retraining a Novelty Detector with Impostor Patterns for Keystroke Dynamics-Based Authentication. In *Advances in Biometrics*, volume 3832 of *Lecture Notes in Computer Science*, pages 633–639. Springer Berlin Heidelberg, 2005.

TABLE XIII

RBF KERNEL PARAMETER ( $\gamma$ ) FOR THE PASCALVOC DATASET.

Aeroplane	Bicycle	Bird	Boat	Bottle
9.5367e-07	9.5367e-07	9.5367e-07	9.5367e-07	3.8147e-06
Bus	Car	Cat	Chair	Cow
2.3842e-07	9.5367e-07	2.3842e-07	1.1921e-07	9.5367e-07
Diningtable	Dog	Horse	Motorbike	Person
4.7684e-07	1.1921e-07	9.5367e-07	4.7684e-07	1.1921e-07
Pottedplant	Sheep	Sofa	Train	Tvmonitor
1.9073e-06	1.9073e-06	2.3842e-07	9.5367e-07	4.7684e-07

- [15] M. M. Moya and D. R. Hush. Network constraints and multi-objective optimization for one-class classification. *Neural Networks*, 9(3):463 – 474, 1996.
- [16] J. Platt. Sequential Minimal Optimization: A Fast Algorithm for Training Support Vector Machines. Technical Report MSR-TR-98-14, Microsoft Research, April 1998.
- [17] D. M. Powers. Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation, 2011.
- [18] G. Ritter and M. T. Gallegos. Outliers in statistical pattern recognition and an application to automatic chromosome classification. *Pattern Recognition Letters*, 18(6):525–539, 1997.
- [19] W. J. Scheirer, L. P. Jain, and T. E. Boult. Probability Models for Open Set Recognition. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 36, November 2014.
- [20] W. J. Scheirer, A. Rocha, A. Sapkota, and T. E. Boult. Towards Open Set Recognition. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 36, July 2013.
- [21] B. Schölkopf and A. J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2001.
- [22] B. Schölkopf, R. C. Williamson, A. J. Smola, J. Shawe-Taylor, and J. C. Platt. Support Vector Method for Novelty Detection. In *Advances in Neural Information Processing Systems 12*, pages 582–588. MIT Press, 2000.
- [23] D. M. Tax and R. P. Duin. Support vector data description. *Machine learning*, 54(1):45–66, 2004.
- [24] A. Vedaldi and A. Zisserman. Efficient additive kernels via explicit feature maps. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 34(3):480–492, 2012.
- [25] T. Yamuna and S. Maheswari. Detection of abnormalities in retinal images. In *Intl. Conf. on Emerging Trends in Computing, Communication and Nanotechnology*, 2013.